# Safe Reinforcement Learning Using Advantage-Based Intervention

Nolan Wagener, Byron Boots, Ching-An Cheng

Georgia Tech · UNIVERSITY of WASHINGTON · Microsoft Research

## Safe RL Problem

Reward $r(s, a) \geq 0 \rightarrow$ Value $V^\pi(s)$

Cost $c(s, a) = \mathbf{1}\{s = \text{violation}\} \rightarrow$ Value $\overline{V}^\pi(s)$

**Goal**: Maximize return while keeping cost below some threshold, *including during training*.

$$\max_\pi \ V^\pi(s_0) \quad \text{subject to} \quad \overline{V}^\pi(s_0) \leq \delta$$

**Dilemma**: Partially optimized RL policy $\pi$ may be unsafe.

**Assumption**: Given baseline policy $\mu$ which is safe starting from the initial state $s_0$.

**Solution**: Use $\mu$ to prevent unsafe actions from $\pi$.

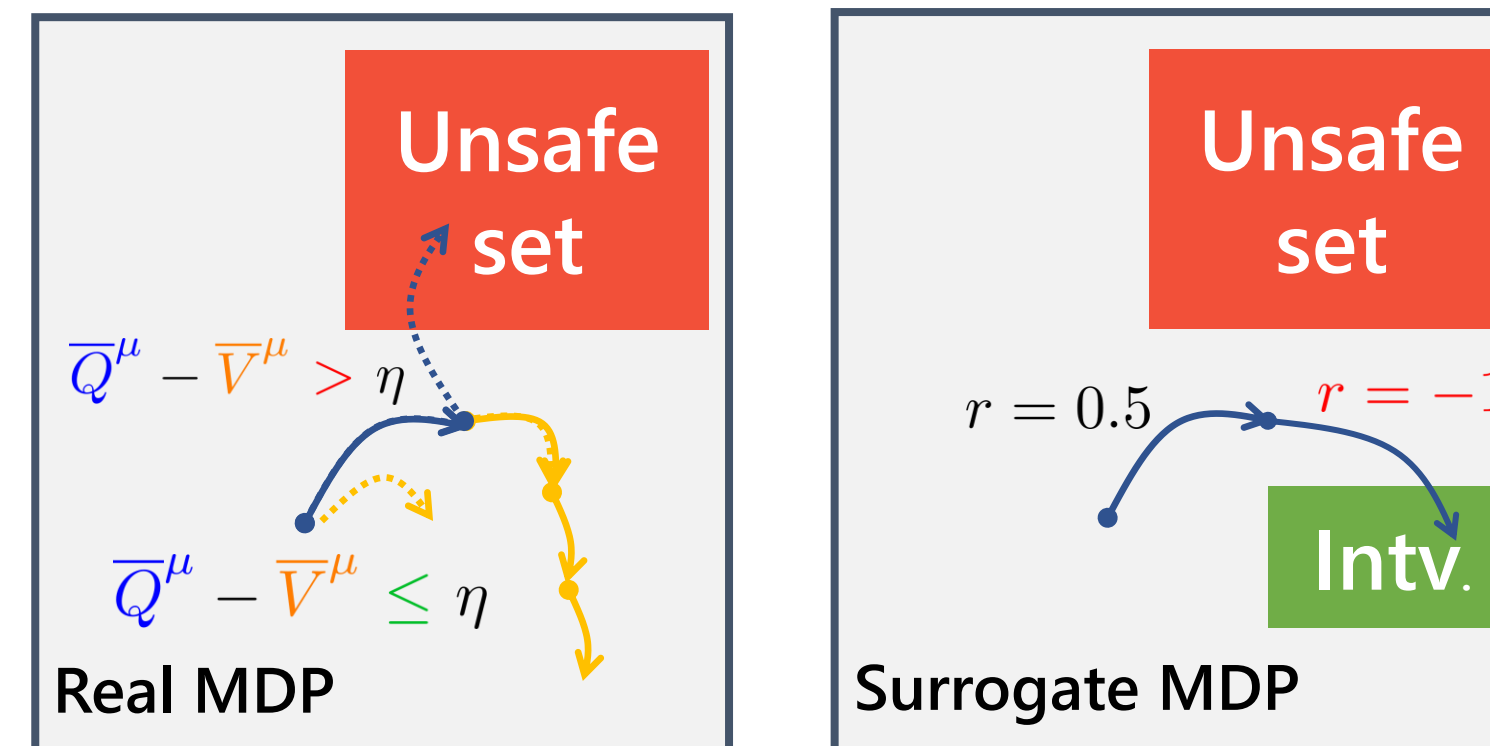**Dilemma**: Optimized $\pi$ needs to be safe and have high returns.

**Solution**: Augment original MDP with penalizing rewards for being intervened by $\mu$.

## Advantage-Based Intervention

Intervention rule $\mathcal{G}$ given by $\mu$ and threshold $\eta$.

Shielded policy $\mathcal{G}(\pi)$ uses advantage $\overline{A}^\mu(s, a)$ w.r.t. $\mu$ to determine to sample from $\pi$ or $\mu$.



What if we don't have access to $\overline{A}^\mu(s, a)$?

Can learn approximation from data collected from $\mu$.

Why intervene based on advantages instead of Q?

Allows reduction from constrained RL to *unconstrained* RL.

## Safe RL with SAILR

Run **shielded policy** $\mathcal{G}(\pi)$ in real MDP.

From $\pi$'s perspective, if intervened it transitions to absorbing state and gets a penalizing reward.

Run regular RL algorithm (e.g., PPO) on surrogate MDP.



Real MDP

Surrogate MDP

## Theoretical Results

**Theorem (Safety During Training)**

Shielded policy $\mathcal{G}(\pi)$ is nearly as safe as $\mu$.

$$\overline{V}^{\mathcal{G}(\pi)}(s_0) \leq \overline{V}^\mu(s_0) + \frac{\eta}{1 - \gamma}$$

**Theorem (Safety and Returns at Deployment)**

SAILR policy $\hat{\pi}$ is nearly as safe as $\mu$ and has nearly the same returns as comparator policy $\pi^*$.

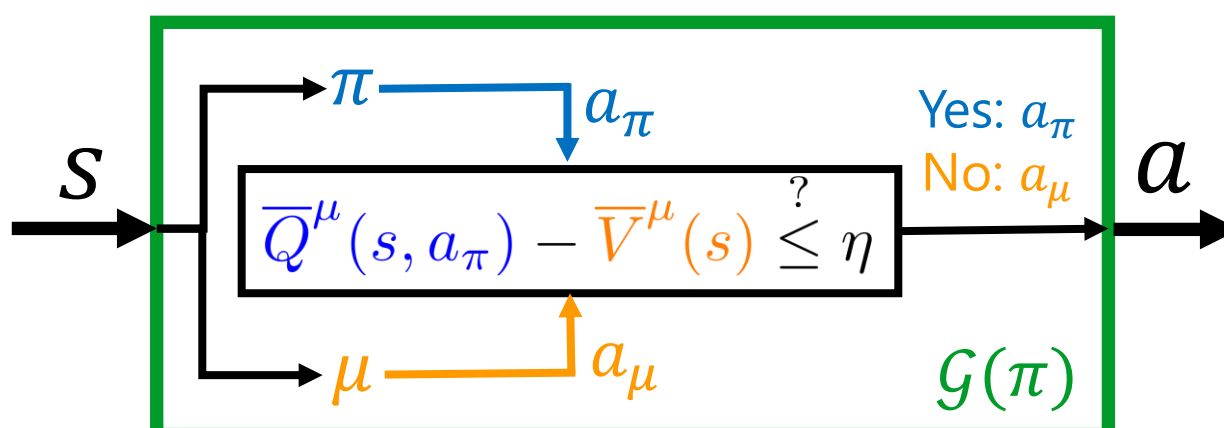$$\overline{V}^{\hat{\pi}}(s_0) \leq \overline{V}^\mu(s_0) + \frac{\eta}{1 - \gamma}$$

$$V^{\pi^*}(s_0) - V^{\hat{\pi}}(s_0) \leq O\left(\frac{\text{Prob}(\pi^* \text{ is intervened by } \mathcal{G})}{1 - \gamma}\right)$$
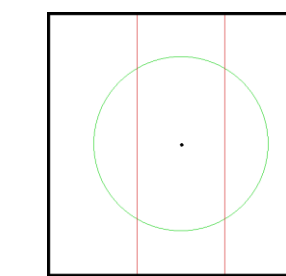
## Experiments

### Point Robot

Reward: move fast in CCW

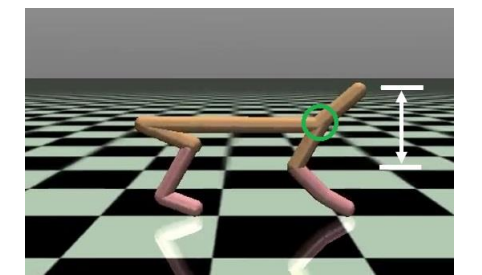Constraint: don't touch vertical red lines
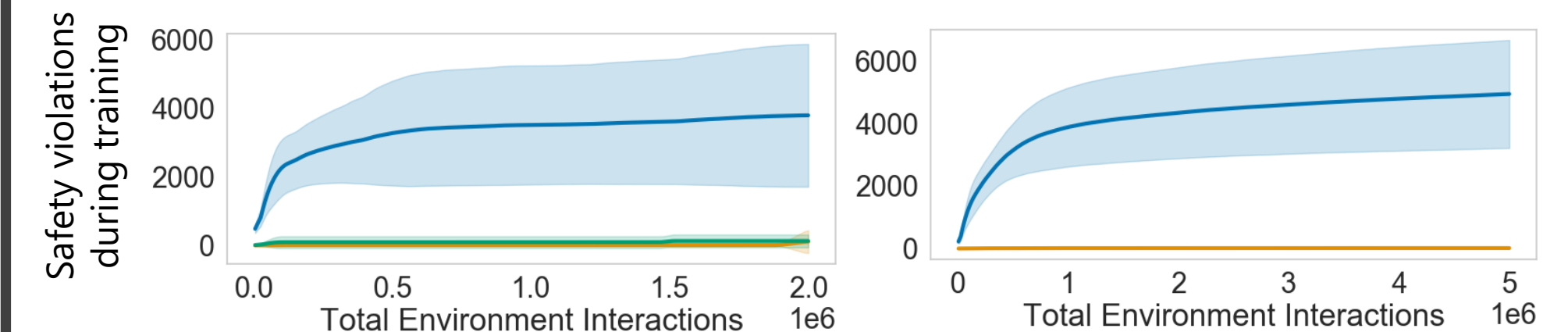
$\mu$: deceleration policy

### Half-Cheetah

Reward: run forward fast
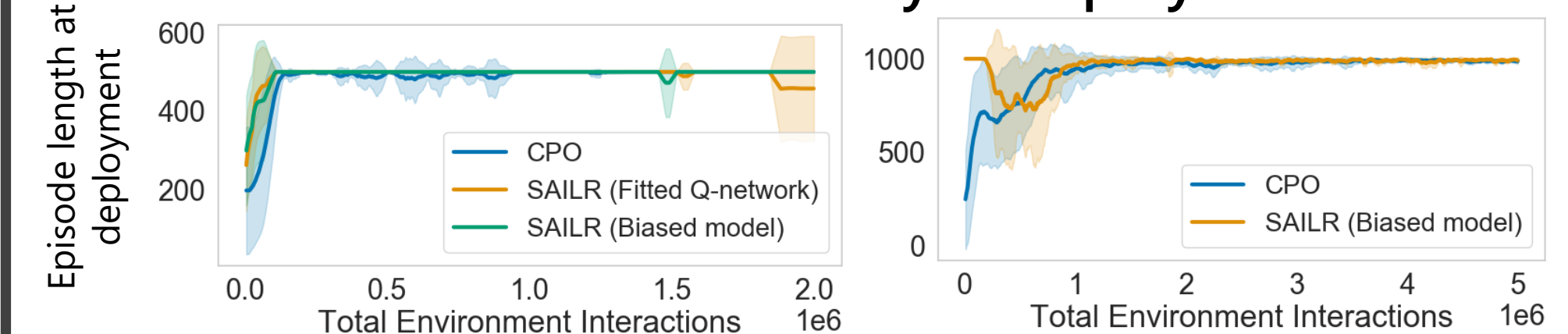
Constraint: keep "chin" joint in a height range

$\mu$: model predictive control



### Far fewer safety violations during training



### Similar level of safety at deployment



### Similar returns at deployment